

Annotations *Rhapsodie*

pour le *Trameur*

13/12/2013 12:51:24

Serge Fleury

Références

Le Trameur, manuel d'utilisation

<http://www.tal.univ-paris3.fr/trameur/leMetierLexicométrique.pdf>

Dans cette documentation, la partie « *Relations de dépendance entre les items de Trame (via leurs annotations)* » présente l'état des développements actuellement disponibles dans le *Trameur* pour l'exploitation des relations de dépendance.

Le Trameur. Propositions de description et d'implémentation des objets textométriques

<http://www.tal.univ-paris3.fr/trameur/trameur-propositions-definitions-objets-textometriques.pdf>

Ce document met au jour une description des objets textométriques et les méthodes mises en œuvre dans le **Trameur** pour travailler sur et avec ces objets dans une perspective textométrique. On y détaille aussi les opérations permises sur une *base textométrique* : format des données textuelles, modification dynamique de la *Trame*, correction ou ajout d'annotation etc.

Sommaire

1. Préambule.....	3
2. Données Rhapsodie	3
3. Intégration des annotations Rhapsodie dans une base textométrique	5
3.1 La Trame textométrique	6
3.2 Le Cadre textométrique.....	8
3.3 Sections.....	10
4. Explorer les relations de dépendance	11
4.1 Recherche de dépendance sur l'ensemble de la base (avec filtrage sur les items en relation)	11
Exemple n°1 : recherche des « objets » du lemme « affirmer »	11
Exemple n°2 : recherche des « sujets » et « objets » du lemme « penser ».....	12
4.2 Retour en contexte	13
4.3 Recherche de dépendance en contexte	14
4.4 Rechercher dans un graphe de dépendance	16
5. Recherche de collocation : spécificités sur relation	19

1. Préambule

Ce document commence par décrire le processus de transcodage des données issues du projet *Rhapsodie* (<http://projet-rhapsodie.fr/>) sous la forme d'une base textométrique importable dans le *Trameur*.

Il présente ensuite les différentes fonctionnalités mises en œuvre pour traiter ce type de données via le *Trameur*.

2. Données Rhapsodie

Les données traitées sont disponibles sur la page du projet *Rhapsodie* :

- Téléchargement des fichiers de codage microsyntactique version bêta 10/13 ([zip](#))
- Tutoriel codage microsyntactique ([pdf](#))

Le fichier d'annotations (*Rhapsodie.tok*) a l'allure suivante (lecture ici dans un tableur) :

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
	TextID	TreeID	TokenID	Token	Lemma	POS	Mode	Tense	Person	Number	Gender	Gov_rection	Type_rection	Gov_para	Type_para	Gov_inher	Type_inher	Gov_junc	Type_junc	Gov_junc_inher	Type_junc_inher
1	A2000	1	1	bonjour	bonjour	0_J						0 root									
2	A2000	1	2																		
3	A2000	1	3	Eric	Eric	0_N				sg	masc	0 root									
4	A2000	2	1	bonjour	bonjour	0_N				sg	masc	0 root									
5	A2000	2	2																		
6	A2000	2	3	À	À	0_Pre						1 dep									
7	A2000	2	4																		
8	A2000	2	5	tout	tout	0_Prs				pl	masc	3 dep									
9	A2000	3	1	nouvelle	nouveau	0_Adj				sg	fem	3 dep									
10	A2000	3	2																		
11	A2000	3	3	mait	mait	0_N				sg	fem	0 root									
12	A2000	3	4																		
13	A2000	3	5	de	de	0_Pre						3 dep									
14	A2000	3	6																		
15	A2000	3	7	pillage	pillage	0_N				sg	masc	5 dep									
16	A2000	3	8																		
17	A2000	3	9	et	et	0_J															
18	A2000	3	10																		
19	A2000	3	11	d	de	0_Pre								5 para_coord		3 dep_inherited		9 junc			
20	A2000	3	12	"	"	1_Pre															
21	A2000	3	13	affrontement	affrontement	0_N				sg	masc	11 dep									
22	A2000	3	14																		
23	A2000	3	15	en	en	0_Pre						3 dep									
24	A2000	3	16																		
25	A2000	3	17	Guadeloupe	Guadeloupe	0_N				sg	masc/fem	15 dep									
26	A2000	4	1	trois	trois	0_D				sg	masc	3 dep									
27	A2000	4	2																		
28	A2000	4	3	policiers	policier	0_N				pl	masc	0 root									
29	A2000	4	4																		
30	A2000	4	5	Messia	Messia	0_Adj				pl	masc	3 dep									

Ces données sont constitués par un certain nombre de textes (l'identifiant du texte est visible dans la première colonne), chacun d'eux est segmenté en « unité illocutoire » (UI, seconde colonne), chacune d'elle est segmentée en *token* (troisième colonne), chacun d'eux est annoté (les autres colonnes)

Le fichier tabulé précédent est composé de 21 colonnes (description à compléter cf *Rhapsodie*) :

1	TextID	Identifiant de l'échantillon (nom de la PARTIE dans la terminologie textométrique)
2	TreeID	Numéro de l'unité illocutoire (UI) dans l'échantillon
3	TokenID	Identifiant du <i>token</i> dans l'UI Les UI d'un échantillon sont séparées les unes des autres par des lignes sans aucun identifiant <i>TreeID</i>
4	Token	Segment de la transcription orthographique pris en 2 blancs ou un blanc et une signe de ponctuation
5	Lemma	Les lemmes sont comme il est d'usage la forme pour les lexèmes invariables, la forme infinitive pour les verbes, le singulier pour les noms et le masculin singulier pour les adjectifs.
6	POS	Partie du discours <ul style="list-style-type: none"> - V pour les verbes - N pour les noms - Adj pour les adjectifs - Adv pour les adverbes - Pre pour les prépositions - CS pour les conjonctions de subordination - J pour les joncteurs : il s'agit des traditionnelles conjonctions de coordinations et d'autres éléments qui lient les couches d'un entassement, comme <i>c'est-à-dire</i> ou <i>y compris</i>. Les éléments clôtureurs d'entassement comme <i>et caetera</i> sont classés comme joncteurs également. - D pour les déterminants - I pour les interjections, y compris des marqueurs de discours comme <i>bon, ben, euh, hein ...</i> - Qu pour les mots qu- que sont les relatifs et les interrogatifs - Cl pour les clitiques, y compris les clitiques sujets (<i>je, tu, il, on, ce</i>) et l'adverbe de négation <i>ne</i>. - Pro pour les autres pronoms - X pour les éléments dont on ne peut déterminer la catégorie syntaxique : partie inaudible (XXX), certaines amorces (quand on ne peut pas deviner le lexème et sa partie du discours), ainsi que les positions non instanciées marquées par &.
7	Mode	Les V reçoivent un trait de mode qui peut prendre 6 valeurs : <i>indicative, subjunctive, imperative, infinitive, past_participle, present_participle</i>
8	Tense	Seuls les V à l'indicatif varient en temps ; le trait <i>tense</i> possède 5 valeurs : <i>present, imperfect, future, conditional</i> et <i>perfect</i>
9	Person	Les V reçoivent aussi des traits d'accord : le trait <i>person</i> a trois valeurs 1, 2 et 3
10	Number	le trait <i>number</i> a deux valeurs <i>sg</i> et <i>pl</i>
11	Gender	le trait <i>genre</i> a deux valeurs <i>fem</i> et <i>masc</i>
12	Gov_rection	
13	Type_rection	
14	Gov_para	
15	Type_para	
16	Gov_inher	
17	Type_inher	
18	Gov_junc	
19	Type_junc	
20	Gov_junc-inher	
21	Type_junc-inher	

Les 11 premières annotations sont réutilisées telles quelles par le processus de transcodage.

Les suivantes sont réutilisées 2 à 2 (cf jeu de couleur) pour construire respectivement une seule annotation (de type relation) qui est réécrite par exemple sous la forme : `Type_rection(Gov_rection)` pour les lignes 12 et 13.

3. Intégration des annotations *Rhapsodie* dans une base textométrique

Les données issues de *Rhapsodie* transcodées dans un format compatible avec le *Trameur* respectent la structuration d'une base textométrique (cf documentation *Trameur*). Celle-ci est composée de 2 parties permettant :

1. La description d'une *Trame* textométrique : liste des items numérotés et annotés (ici chaque item est associé à 13 annotations)
2. La description du *Cadre* textométrique : liste des partitions définies sur la trame ; chacune porte un nom et est associée à une liste de parties définies chacune par son nom (le nom de l'échantillon de *Rhapsodie*), par sa position de début sur la *Trame* et sa position de fin

Le fichier issu du transcodage est au final une base textométrique qui se charge dans le *Trameur* par le module « *importation de base* ».

nous ferons le point sur les violences en direct avec notre envoyé spécial Sébastien Paour dans quelques instants \$

Position: <91>
 Forme: <point>|Freq: 34
 Lemme: <point>|Freq: 23
 Cat: <B_N>|Freq: 6247
 a-00004: <->|Freq: 70852
 a-00005: <->|Freq: 72497
 a-00006: <->|Freq: 68133
 a-00007: <sg>|Freq: 17353
 a-00008: <masc>|Freq: 8816
 a-00009: <OBJ(87)>|Freq: 1
 a-00010: <->|Freq: 74709
 a-00011: <->|Freq: 74145
 a-00012: <->|Freq: 75936
 a-00013: <->|Freq: 76812

n°Annotation	Label	Contenu
1	Forme	Forme graphique
2	Lemme	Lemme
3	Cat	P.O.S
4	a-00004	Mode
5	a-00005	Tense
6	a-00006	Person
7	a-00007	Number
8	a-00008	Gender
9	a-00009	Type_rection(Gov_rection)
10	a-00010	Type_para(Gov_para)
11	a-00011	Type_inher(Gov_inher)
12	a-00012	Type_junc(Gov_junc)
13	a-00013	Type_junc-inher(Gov_junc-inher)

3.2 Le Cadre textométrique

Les différents échantillons initiaux de *Rhapsodie* sont considérés comme autant de parties différentes : la base finale est donc une partition de textes (*Cadre*), chaque partie contient les zones textuelles associées à l'identifiant initial de l'échantillon. Ci-dessous, le *Cadre* final mis au jour dans le *Trameur* et son codage dans la base construite après transcodage :

3.3 Sections

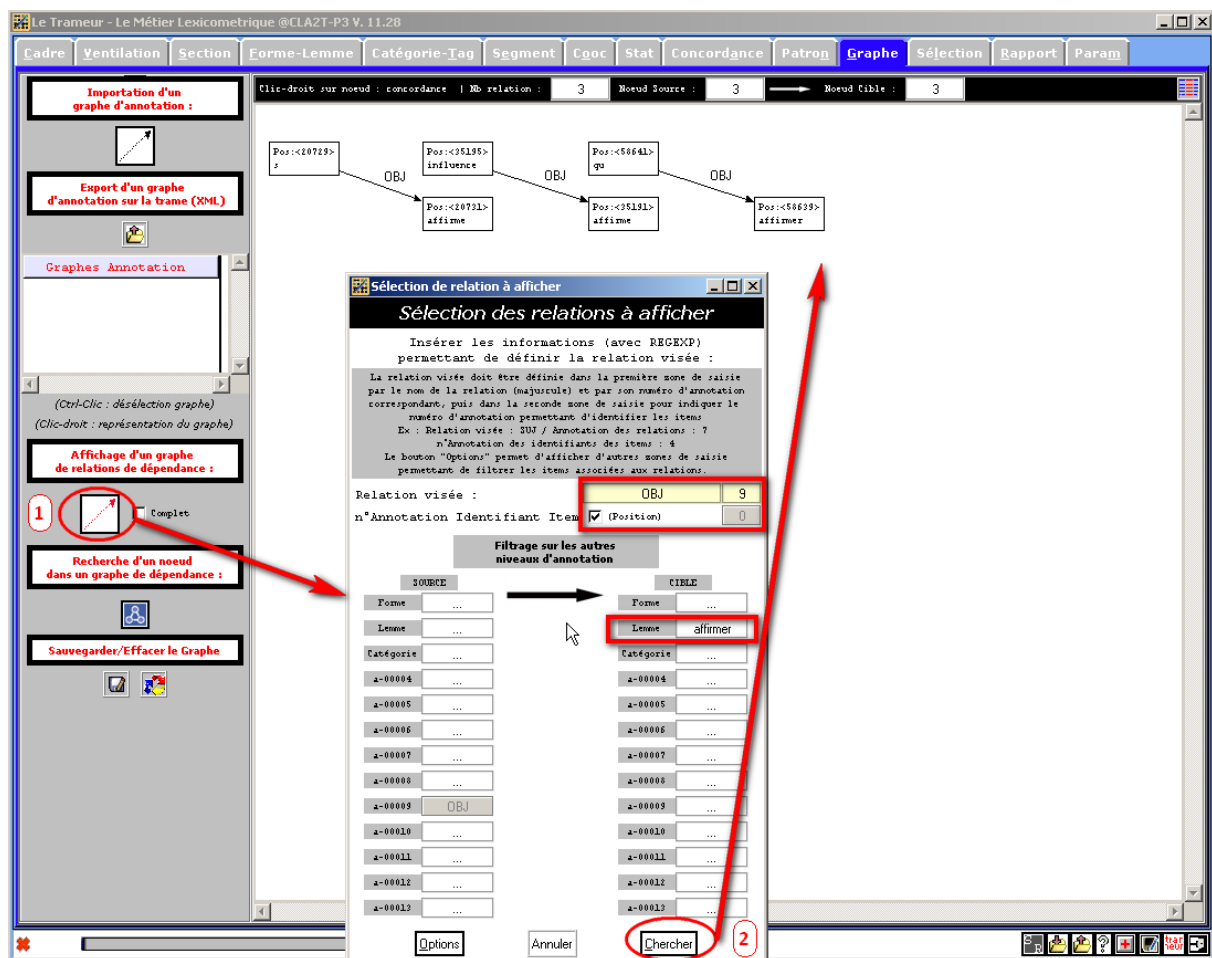
Le processus de transcodage intègre aussi un marquage de sections : après chaque UI, un caractère délimiteur de section (§) est introduit (*cf* caractère en position 5 dans la base présentée ci-dessus) pour permettre de construire dans le *Trameur* une représentation cartographique de la base sous la forme d'une *carte des sections* :

4. Explorer les relations de dépendance

Les différentes fonctionnalités disponibles dans le *Trameur* pour travailler avec les annotations de relations sont décrites dans la documentation du *Trameur*. On les illustre ci-dessous sur les données de la base *Rhapsodie2Trameur*.

4.1 Recherche de dépendance sur l'ensemble de la base (avec filtrage sur les items en relation)

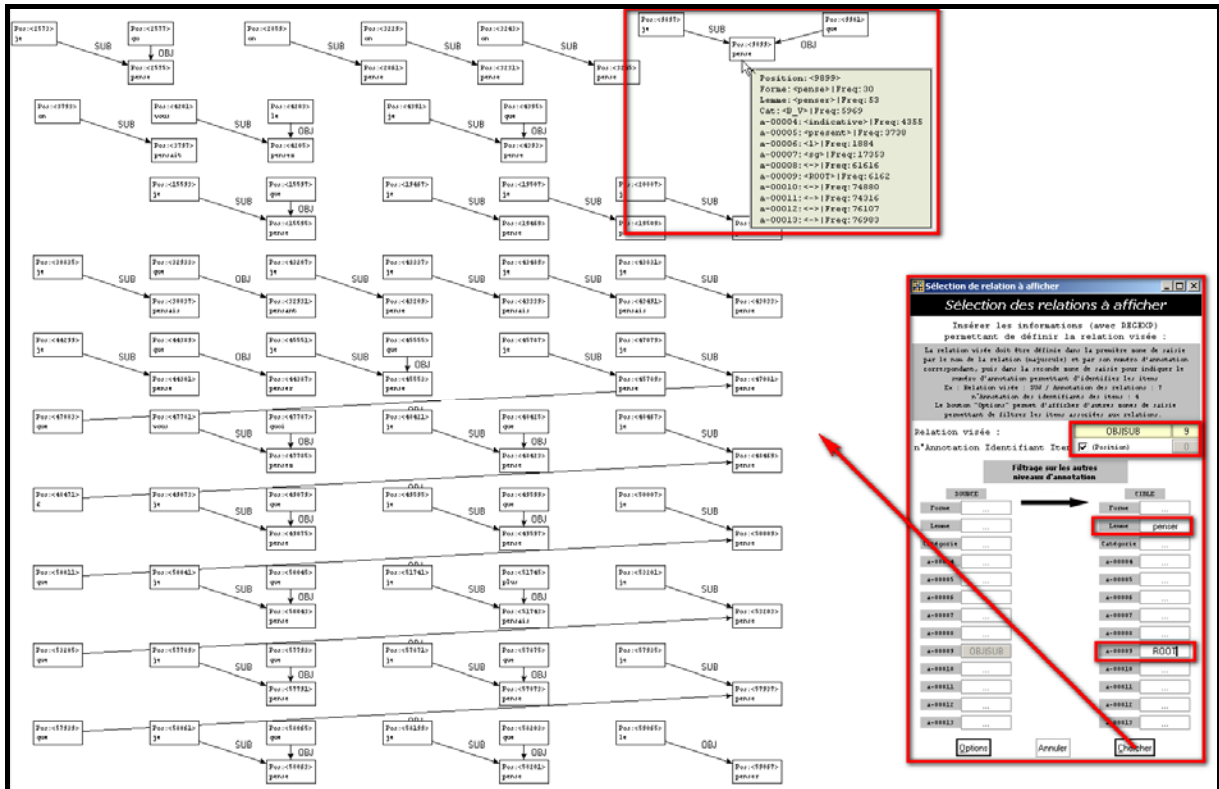
Exemple n°1 : recherche des « objets » du lemme « affirmer »



Dans la figure précédente, on cherche des relations de dépendance de type OBJ en contraignant un des composants de la relation : il doit porter le lemme "affirmer" (i.e. on cherche les objets d'affirmer).

Exemple n°2 : recherche des « sujets » et « objets » du lemme « penser »

Dans la figure suivante, la relation cherchée est double via l'expression régulière SUB|OBJ i.e SUB ou OBJ, la cible de la relation impose une valeur pour le lemme (« penser ») et pour l'annotation n°9 (ROOT)

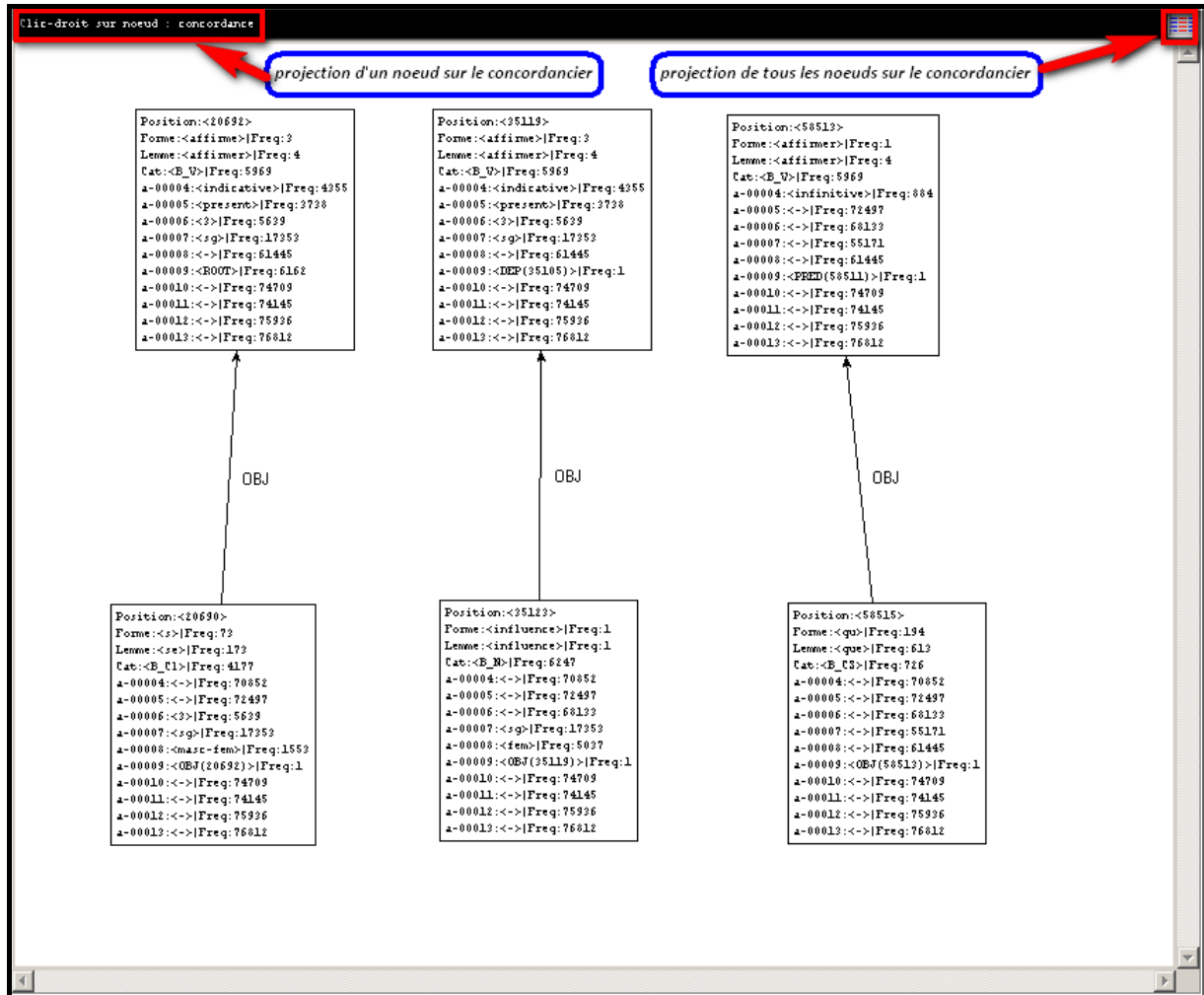


Dans les 2 graphes précédents, l'affichage des nœuds est celui produit par défaut : affichage de la position de l'item sur la *Trame* et de sa forme graphique. On peut visualiser l'ensemble des annotations d'un nœud donné via le mode Aperçu disponible dans tous les éditeurs du Trameur (en passant la souris au-dessus de l'objet visé) : ci-dessus, l'item en position 9899 (lemme : penser) est mis au jour.

Dans la figure qui suit, l'affichage d'un nœud est directement rendu « complet » via l'option du même nom disponible et activable avant de charger un graphe.

4.2 Retour en contexte

Un graphe de relations de dépendance étant produit, chaque nœud du graphe peut-être projeté sur le concordancier (clic-droit sur le nœud). L'ensemble des nœuds peut aussi être projeté globalement sur le concordancier :



La concordance produite dans ce dernier cas a l'allure suivante :

```

-----PARTIE{PARTIE=M2004}-----
      § un nouvel ordre international s'affirme peu à peu
-----PARTIE{PARTIE=M2001}-----
      qu'une nation affirme son influence qu'elle pèse dans
-----PARTIE{PARTIE=D2013}-----
      Rachida Dati pour oser affirmer qu'il ne s'agit
  
```

Les sources de la relation sont coloriées en rouge (pôle de la concordance), les cibles en vert (*i.e* ici le verbe est en vert, son objet en rouge).

On peut aussi varier l'affichage des annotations dans la concordance en matérialisant par exemple la catégorie de chaque item :

```

PARTIE{PARTIE=M2004}
$ un-B_D nouvel-B_Adj ordre-B_N international-B_Adj s-B_C1 '-I_C1 affirme-B_V peu-B_Adv à-I_Adv peu-I_Adv
PARTIE{PARTIE=M2001}
qu-B_Qu '-I_Qu une-B_D nation-B_N affirme-B_V son-B_D influence-B_N qu-B_Qu '-I_Qu elle-B_C1 pèse-B_V dans-B_Pre
PARTIE{PARTIE=D2013}
Rachida-B_N Dati-B_N pour-B_Pre oser-B_V affirmer-B_V qu-B_CS '-I_CS il-B_C1 ne-B_C1 s-B_C1 '-I_C1 agit-B_V

```

Toutes les zones d'édition du *Trameur* permettent de mettre au jour les annotations de la *Trame* (cf documentation en ligne, partie « Marquage des annotations de la *Trame* »). Si on considère la figure suivante, elle présente un extrait de la concordance construite à partir du graphe ayant permis d'extraire la relation OBJ (cf graphe infra) :

```

PARTIE{PARTIE=H002}
$ de la crise aux Antilles il ne devrait pas en être question au sommet social de l'Élysée cet
savamment minuitée $ alors avant d'aborder ce qui peut en sortir je vous propose de voir où quand et
alors avant d'aborder ce qui peut en sortir je vous propose de voir où quand et comment cela va
d'aborder ce qui peut en sortir je vous propose de voir où quand et comment cela va se passer
passer avec vous Jean-François Achilli $ alors passons maintenant au détail des mesures discutées et aux attentes des syndicats
des mesures de justice je cite pour les salariés touches par la crise économique $ mais Sara Ghibaudo il ne
vous d'aujourd'hui les Français se montrent très sévères à l'égard de la politique économique du gouvernement qu'
des voitures $ bonjour Sébastien Paour $ vous vous trouvez au Gosier $ où en est la situation $ merci
$ tout le monde est maintenant en tout cas suspendu à ce que Nicolas Sarkozy annoncera demain lors de sa
a promulgué son plan de relance adopté ce week-end par le congrès $ sept cent quatre-vingt-sept milliards
$ sept cent quatre-vingt-sept milliards de dollars destinées à sauver ou créer plus de trois millions et demi
à ralentir les saisies immobilières devrait être annoncée aujourd'hui par le président américain $ cela concerne plusieurs millions d'
Motors et Chrysler ont présenté hier leur plan de restructuration à l'administration Obama $ ils demandent vingt-deux milliards
toujours en garde à vue $ les enquêteurs cherchent maintenant à savoir si elle a été complice dans cette évasion
a été complice dans cette évasion $ elle se trouvait dans le parloir de la prison au moment où les
la foire d'empoigne $ la tension est encore montée d'un cran hier $ les avocats du berger corse
procès en attendant que la cour d'assise se prononce sur leur demande de supplément d'information $ lundi le
d- ce sont des procédés terroristes $ il le dit à Laurent Doulsant $ direction maintenant Barcelone pour un congrès
Laurent Doulsant $ direction maintenant Barcelone pour un congrès consacré au téléphone mobile en crise lui aussi $ pour la
nul zéro partout pour Lille au Mans $ Valenciennes sort de la zone de relégation après sa victoire sur Caen
sa victoire sur Caen deux zéro $ Nice s'impose à Nancy deux à un $
PARTIE{PARTIE=H002}
voit euh qui voit ça passer $ et donc elle se dit eh bah elle peut récupérer elle pourra récupérer
rencontrent $ et et donc elle elle tombe $ ils se $ $ ouais il y a un accident quoi
ouais il y a un accident quoi $ et ils se $ $ ils tombent tous les deux $ et c'
boulanger voit que la baguette a disparu $ et il se dit que $ $ ouais moi je pense qu'
va $ $ alors le euh euh le boulanger dit au policier qu'en fait c'est la fille $

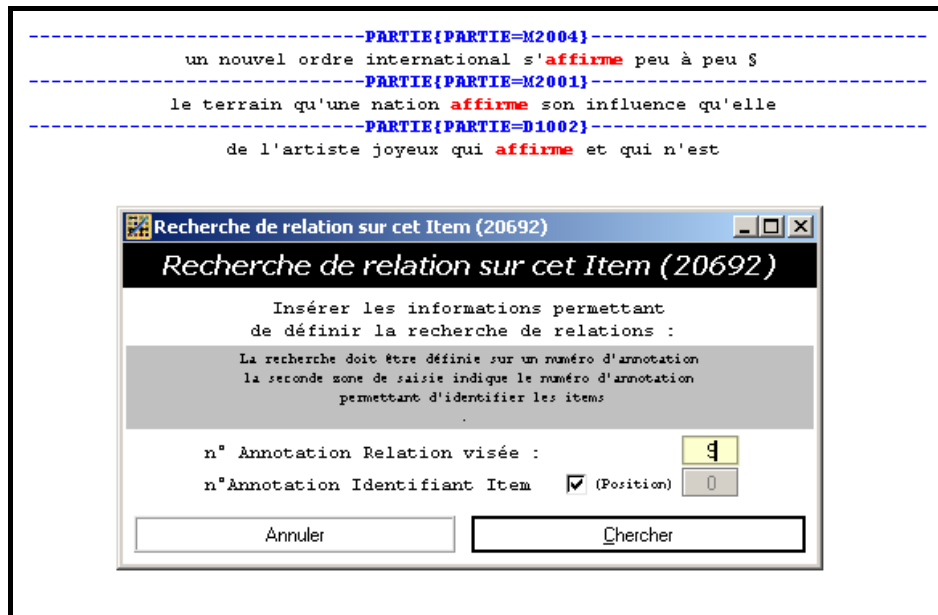
```

Le marquage de certains items est possible en sélectionnant les annotations à mettre au jour :

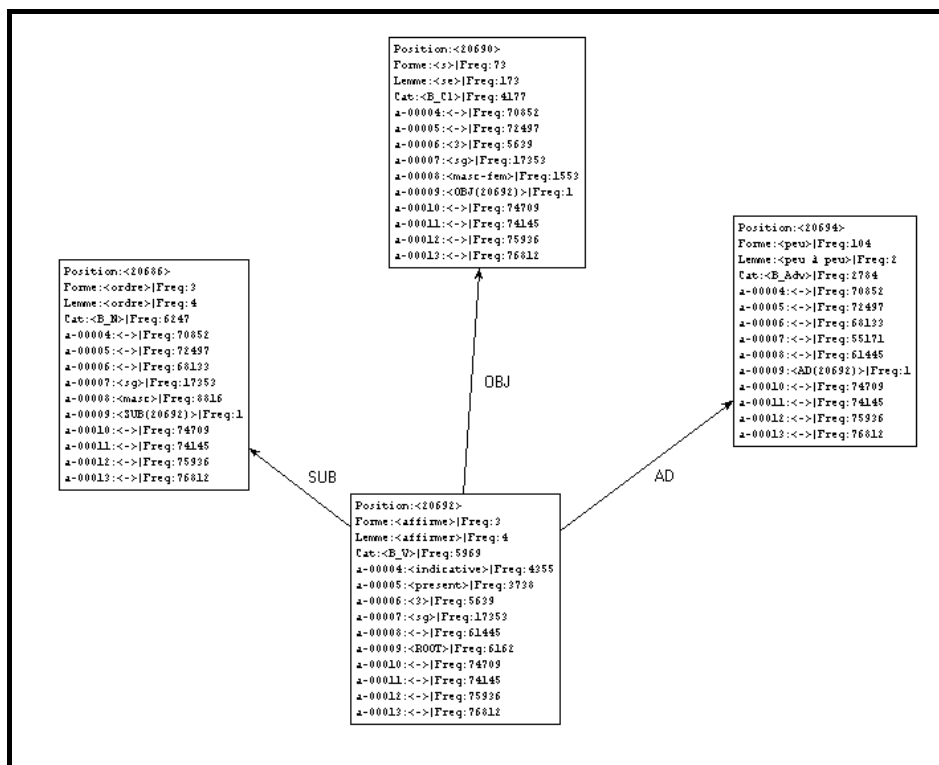
The screenshot shows the Trameur software interface. On the left, there is a sidebar with a search bar and a list of annotations. The main area displays a concordance table with text from the 'Trame' and its annotations. A dialog box titled 'Sélection des annotations à colorier...' is open, showing a list of annotations with checkboxes. A red arrow points from the dialog box to the concordance table, indicating the selection process.

4.3 Recherche de dépendance en contexte

La recherche de dépendance peut aussi être réalisée en contexte, par exemple sur le concordancier. Si on considère la concordance suivante (en haut de la figure) :



Le raccourci clavier Ctrl-Clic-droit sur un item de la concordance (ici l'item visé est la première occurrence de la forme graphique « affirme ») permet de rechercher toutes les relations de dépendance pointant sur cet item. On commence par indiquer où chercher les relations de dépendance (numéro d'annotation portant ce type d'information : ici le n°9) et comment sont indexer les items sur la trame (ici par leur position). La recherche conduit à la production d'un graphe mettant au jour toutes les relations sur l'item visé :



Le graphe donne à voir les 3 relations pointant sur l'item.

4.4 Recherche dans un graphe de dépendance

On présente tout d'abord le résultat produit par la requête suivante « recherche de la relation OBL » :

Sélection de relation à afficher

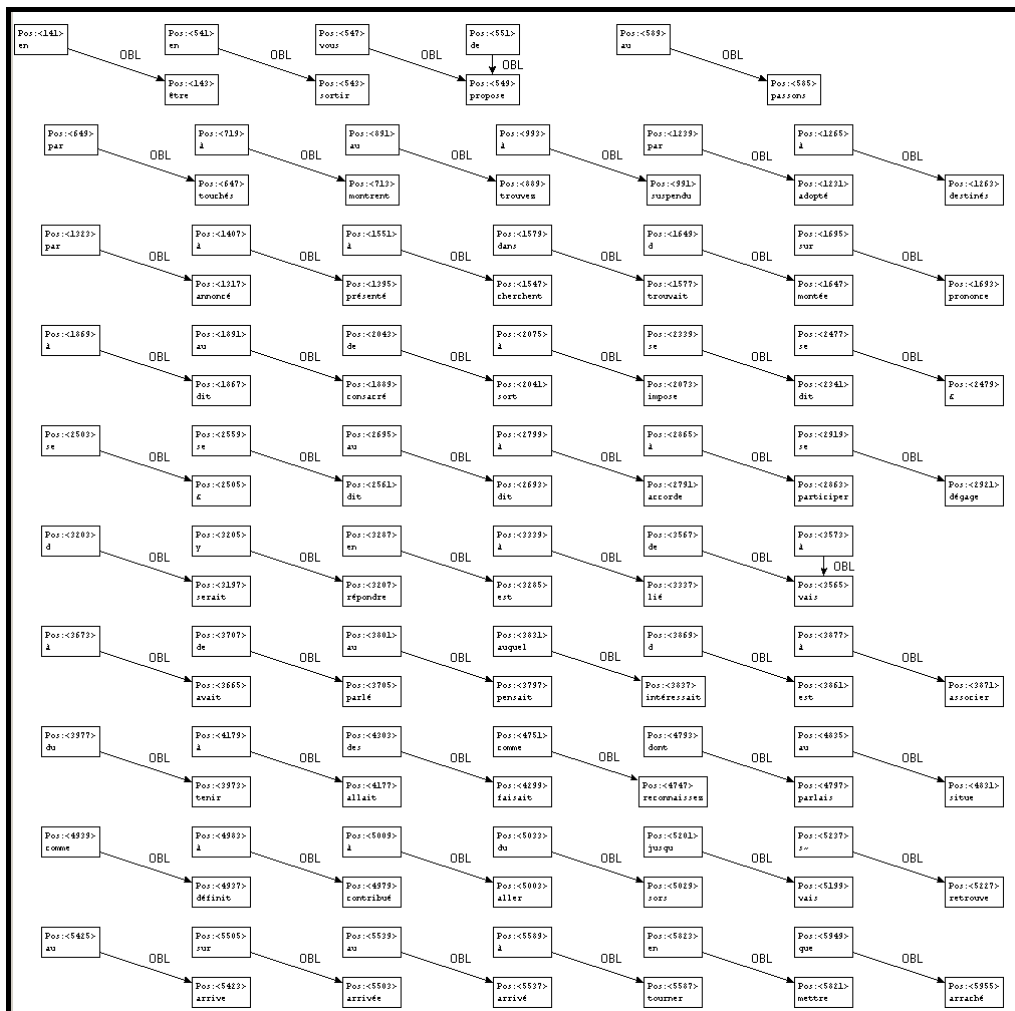
Sélection des relations à afficher

Insérer les informations (avec REGEXP)
 permettant de définir la relation visée :
 La relation visée doit être définie dans la première zone de saisie
 par le nom de la relation (majuscule) et par son numéro d'annotation
 correspondant, puis dans la seconde zone de saisie pour indiquer le
 numéro d'annotation permettant d'identifier les items
 Ex : Relation visée : OBL / Annotation des relations : ?
 n°Annotation des identifiants des items : 4
 Le bouton "Options" permet d'afficher d'autres zones de saisie
 permettant de filtrer les items associées aux relations.

Relation visée :

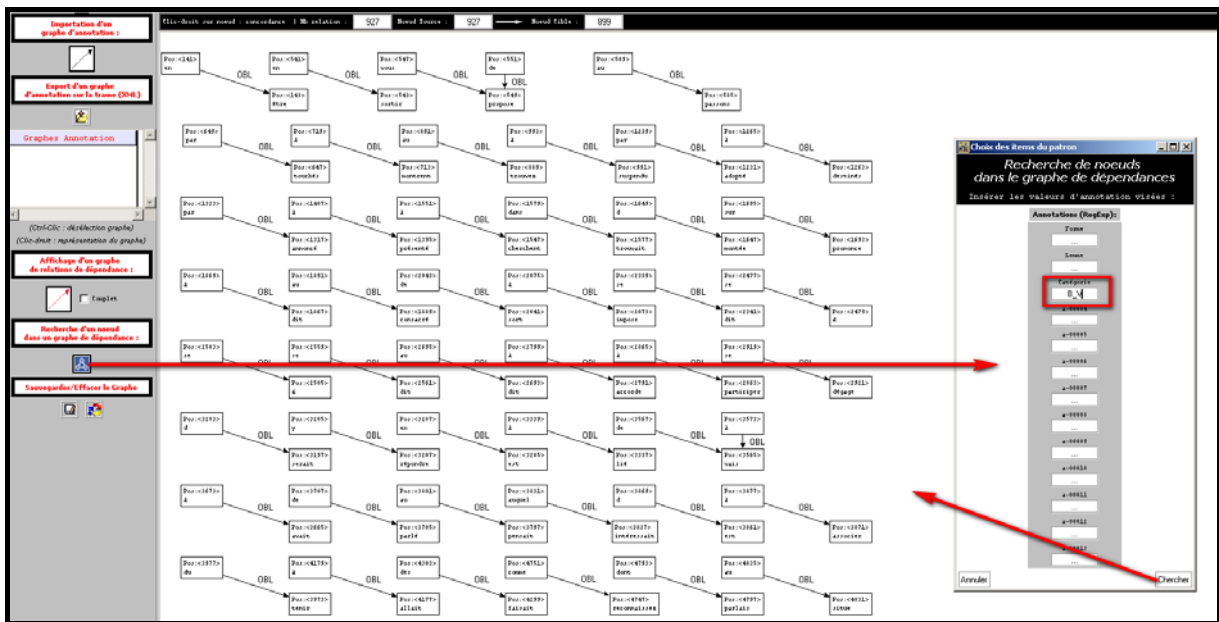
n°Annotation Identifiant Item ☒ (Position)

Le résultat produit a l'allure suivante :

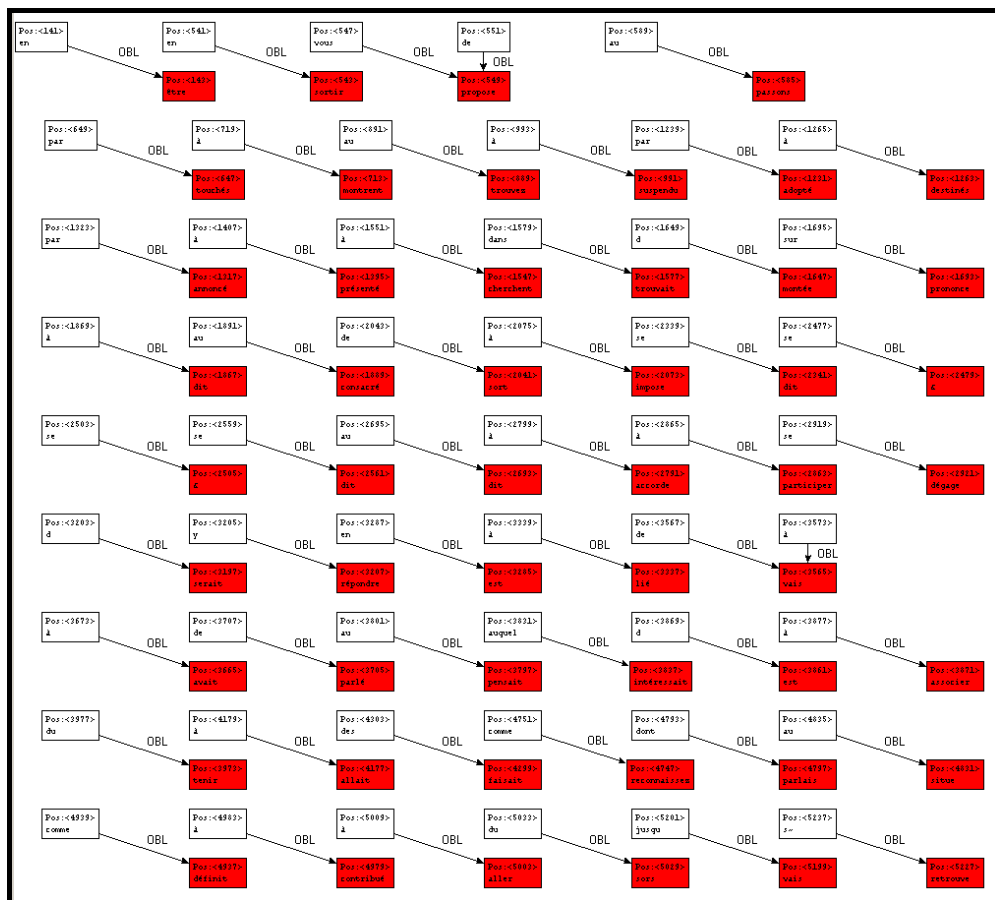


On peut mettre au jour dans ce graphe certains nœuds en filtrant leurs annotations.

Le processus de filtrage des nœuds du graphe permet de sélectionner certains nœuds sur la base des différentes annotations disponibles. Dans l'exemple suivant, on veut mettre au jour les nœuds de catégorie B_V :



Le résultat produit a l'allure suivante :



Dans la figure qui suit, 3 filtrages successifs sont réalisés pour marquer successivement les nœuds ayant pour catégorie B_V, puis B_N et enfin B_Pre; chaque requête est précédée par la modification de la couleur à utiliser pour le marquage des nœuds visés (verbe en bleu, nom en vert et préposition en orange) :



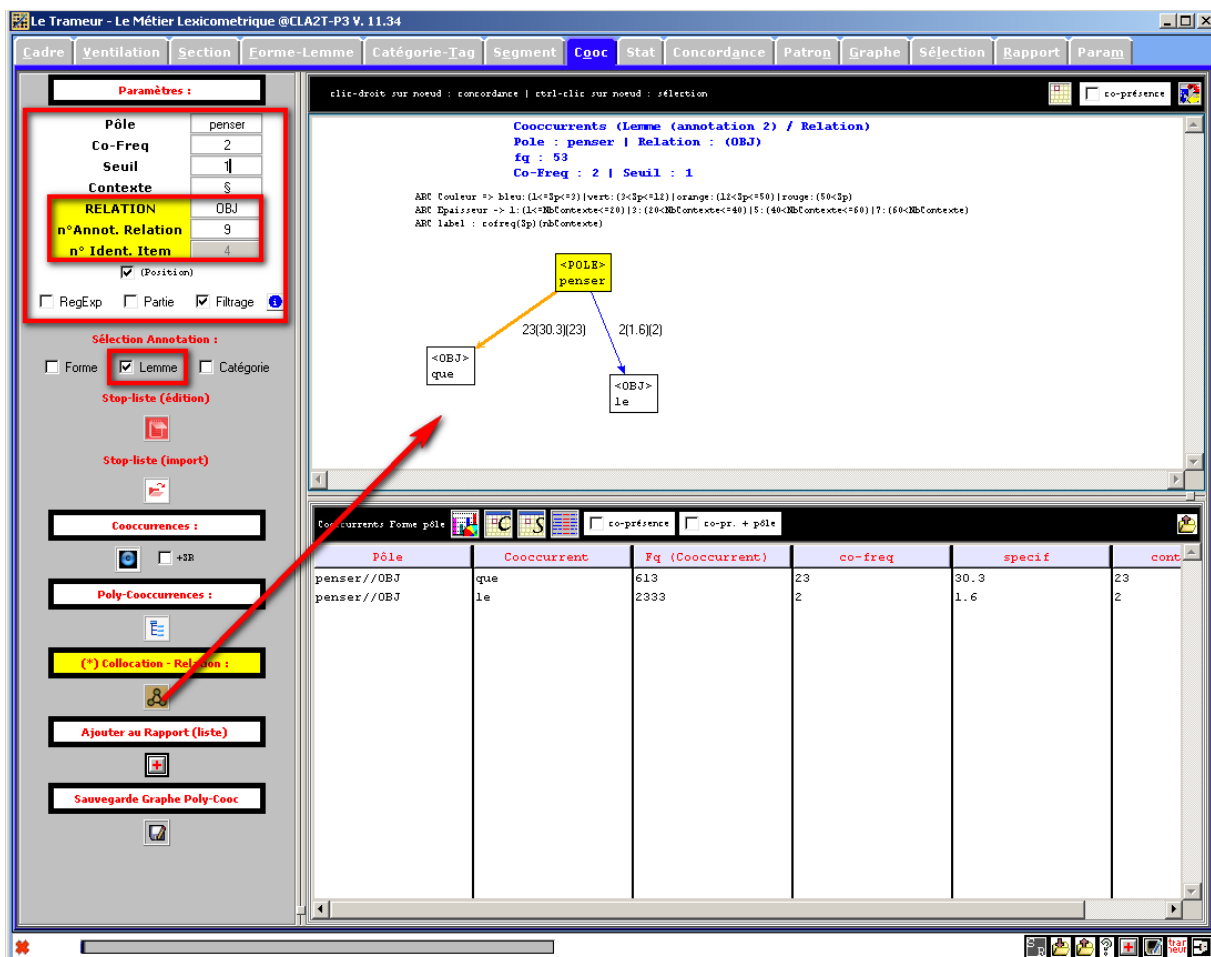
5. Recherche de collocation : spécificités sur relation

Une version particulière du module de calcul des cooccurrences permet de prendre en compte les relations entre les items de la Trame :

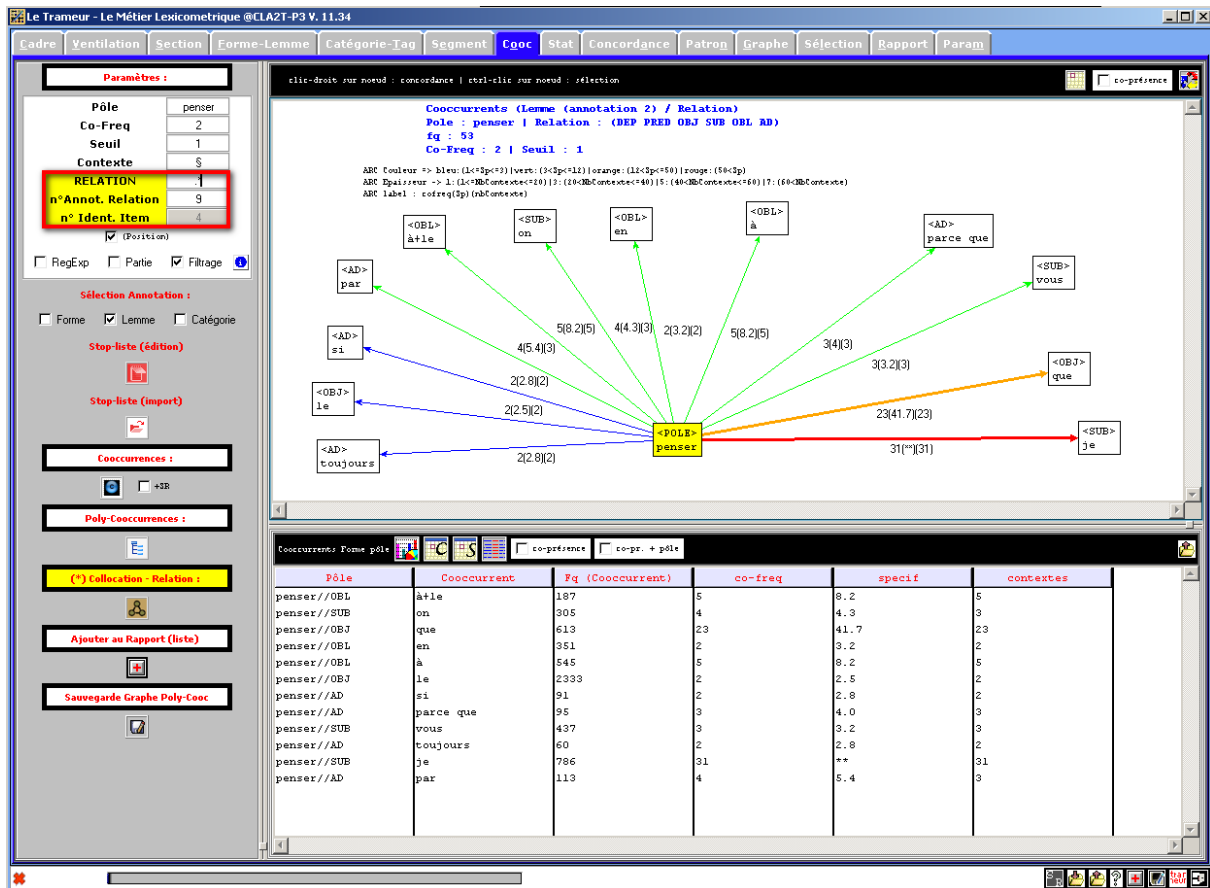
Etant donné une forme pôle, la recherche de ses cooccurents peut être contrainte par la prise en compte d'une relation entre ce pôle et ses candidats cooccurents. Pour un item X donné (le pôle), on s'intéresse aux items Y cooccurents de X et en relation REL avec X (les collocatifs de X) :

X -> REL -> Y

La figure suivante présente les cooccurents du lemme penser en contraignant ses cooccurents à être en position d'objet :



La relation entre le pôle et ses cooccurents peut être « non déterminée » par avance, dans la figure suivante, la relation est exprimée sous la forme : .* (regexp) ; elle vise donc toutes les relations du type : penser-OBJ-y, penser-SUB-y etc.



Une fois le graphe construit, un clic-droit sur un des noeuds montrent les contextes dans lesquels la relation se réalise (*i.e* les contextes utilisés par le calcul).

Dans cet exemple, 3 contextes dans lesquels *vous* est sujet de *penser* (*cf* arc *penser* → *vous*) ont été mis au jour, un clic-droit sur le nœud *vous* les montrent dans le concordancier (les 2 items de la relation y sont colorés automatiquement : *vous* en rouge (pôle de la concordance) et *penser* en vert).

<p>-----PARTIE{PARTIE=D2009 }-----</p> <p>un mot DELIM ou bien vous le penser vraiment DELIM</p> <p>-----PARTIE{PARTIE=D0006 }-----</p> <p>dans le quartier de Paris vous en penser quoi ici</p> <p>-----PARTIE{PARTIE=D2011 }-----</p> <p>DELIM que^êtreUNKNOWNce que vous en penser de le</p>	<p>Position: <64987></p> <p>Forme: <vous> Freq: 441</p> <p>Lemme: <vous> Freq: 437</p> <p>Cat: <B_C1> Freq: 4177</p> <p>a-00004: <-> Freq: 71023</p> <p>a-00005: <-> Freq: 72668</p> <p>a-00006: <2> Freq: 988</p> <p>a-00007: <p1> Freq: 3992</p> <p>a-00008: <-> Freq: 61616</p> <p>a-00009: <SUB(64991)> Freq: 1</p> <p>a-00010: <-> Freq: 74880</p> <p>a-00011: <-> Freq: 74316</p> <p>a-00012: <-> Freq: 76107</p> <p>a-00013: <-> Freq: 76983</p>
---	---